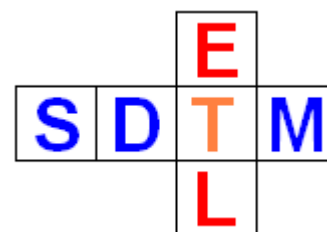


SDTM-ETL 4.3 User Manual and Tutorial

Author: Jozef Aerts, XML4Pharma

Last update: 2023-08-28



Tutorial: Merging datasets

Table of Contents

Introduction.....	1
SDTM-ETL and 'split' datasets	2
Merging datasets that were generated separately.....	3
Conclusions.....	15

Introduction

SDTM and regulatory authorities don't make it easy to us when it comes to generating and submitting comprehensive and especially logical, sets of data.

For example, even now in 2023, it forces us to "ban" non-standard variables (i.e. sponsor-defined variables) to a supplemental qualifier (SUPPxx) dataset. Also, when a data point value exceeds 200 characters, it must be split into chunks of not more than 200 characters and the second and further chunk must be banned to SUPPxx. Also, there is the famous 5GB file size limit, forcing us to "split" datasets¹ when the XPT file size grows beyond that limit.

Essentially, there are three reasons for all this:

The first is the mandated use of the outdated SAS Transport 5 (XPT) format, which is [extremely inefficient in byte storage](#). Selecting XPT 30 years ago was a major error, even simple CSV would have been a better choice.

The second is the lack of SDTM understanding at the FDA: many of the reviewers are not able to distinguish between standard and non-standard SDTM variable, so CDISC decided that these need to go into SUPPxx datasets, which reviewers are supposed to then re-merge into the "parent" dataset. This could easily be solved e.g. by color-coding of non-standard variable columns in modern viewers (as the information about standard and non-standard is in the define.xml), but the relative primitive tools (third reason) at the regulatory authorities do not support this.

Modern viewers such as the "[Smart Submission Dataset Viewer](#)" have such features, e.g.:

¹ The interesting thing is that in such case, one still needs to submit the >5GB dataset, although FDA claims that it cannot read these ...

AGE	AGEU	SEX	RACE	ETHNIC	ARMCD	ARM	ACTARMCD	ACTARM	COUNTRY	DMDC	DMDY	COMPLT16	COMPLT24	COMPLT8	EFFICACY	SAFETY	ITT
61	YEARS	F	WHITE	HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-12-26	-7	Y	Y	Y	Y	Y	Y
64	YEARS	M	WHITE	HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2012-07-22	-14						
71	YEARS	M	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2013-07-11	-8	Y	Y				
74	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2014-03-10	-8						
77	YEARS	F	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2014-06-24	-7	Y	Y				
85	YEARS	F	WHITE	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-01-22	-21						
59	YEARS	F	WHITE	HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2013-12-20							
68	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2013-12-23	-9	Y	Y				
81	YEARS	F	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2012-08-25	-13						
84	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2012-11-23	-7	Y	Y	Y	Y	Y	Y
52	YEARS	M	WHITE	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2014-02-27	-13	Y	Y	Y	Y	Y	Y
84	YEARS	M	WHITE	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2014-02-09	-6	Y	Y	Y	Y	Y	Y
81	YEARS	F	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2012-10-23	-5	Y	Y	Y	Y	Y	Y
57	YEARS	F	WHITE	NOT HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2013-09-05							
75	YEARS	F	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2013-05-07	-13			Y	Y	Y	Y
57	YEARS	M	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2013-08-14	-9	Y	Y	Y	Y	Y	Y
79	YEARS	F	WHITE	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-09-06	-17	Y	Y	Y	Y	Y	Y
82	YEARS	F	WHITE	NOT HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2013-04-18							
62	YEARS	F	AMERICAN...	NOT HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2012-09-30							
56	YEARS	M	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2013-01-28	-15			Y	Y	Y	Y
79	YEARS	F	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Lo	Xanomelin...	USA	2013-11-26	-9			Y	Y	Y	Y
71	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2013-02-03	-12			Y	Y	Y	Y
80	YEARS	F	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2012-07-08	-14	Y	Y	Y	Y	Y	Y
81	YEARS	F	BLACK OR ...	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-01-25	-8	Y	Y	Y	Y	Y	Y
76	YEARS	F	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2012-10-30	-16			Y	Y	Y	Y
69	YEARS	M	WHITE	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-03-20	-10	Y	Y	Y	Y	Y	Y
56	YEARS	M	WHITE	HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2013-12-28	-14	Y	Y	Y	Y	Y	Y
57	YEARS	F	BLACK OR ...	NOT HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2013-09-22							
61	YEARS	M	AMERICAN...	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2014-01-25	-13	Y		Y	Y	Y	Y
56	YEARS	F	WHITE	HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2014-01-17	-8	Y	Y	Y	Y	Y	Y
67	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2013-03-17	-7			Y	Y	Y	Y
61	YEARS	M	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Hi	Xanomelin...	USA	2013-08-20	-9			Y	Y	Y	Y
80	YEARS	F	WHITE	NOT HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2013-12-08							
68	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2014-05-10	-12	Y	Y	Y	Y	Y	Y
79	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2012-09-16	-16	Y	Y	Y	Y	Y	Y
51	YEARS	M	WHITE	NOT HISPANIC	Xan_Lo	Xanomelin...	Xan_Lo	Xanomelin...	USA	2012-12-22	-14			Y	Y	Y	Y
63	YEARS	F	WHITE	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-10-01	-7	Y	Y	Y	Y	Y	Y
54	YEARS	F	WHITE	HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2014-04-01							
67	YEARS	M	WHITE	NOT HISPANIC	Xan_Hi	Xanomelin...	Xan_Lo	Xanomelin...	USA	2013-07-24	-7			Y	Y	Y	Y
81	YEARS	F	BLACK OR ...	NOT HISPANIC	Pbo	Placebo	Pbo	Placebo	USA	2013-05-20	-10	Y	Y	Y	Y	Y	Y
74	YEARS	M	WHITE	NOT HISPANIC	Scrfail	Screen Fail...	Scrfail	Screen Fail...	USA	2013-09-18							

(the DM file is from the famous LZTZ CDISC-FDA SDTM pilot)

SDTM-ETL and 'split' datasets

With SDTM-ETL, we never need to "split" datasets, we take care in advance that we never produce huge datasets (larger than 5GB) that we then later need to split. Even for large datasets smaller than 5GB, we need to think about whether such huge datasets are "usable" for reviewers. For example, an LB dataset with millions of rows spread over different categories of lab tests will be hard to analyze for reviewers - they will need to start to filter to make these usable.

Therefore, the better strategy, followed by SDTM-ETL, is to develop different instances of the same domain, usually one dataset per category. This also makes the mapping considerably easier, as data for different categories can come from different sources, e.g. some from EDC, other from electronic data transfers (e.g. in CSV format).

This strategy e.g. leads to different dataset definitions in SDTM-ETL, for example:

RELREC		STUDYID	RDOMAIN	USUBJID	IDVAR
SUPPQUAL		STUDYID	RDOMAIN	USUBJID	IDVAR
RELSUB		STUDYID	USUBJID	RELSUB.POOLID	RELSUB.RS
OI		STUDYID	DOMAIN	OI.NHOID	OI.OISEQ
██████████	GLOBAL	STARTREF			
██████████	LBUR	STUDYID	DOMAIN	USUBJID	LB.LBSEQ
██████████	LBBL	STUDYID	DOMAIN	USUBJID	LB.LBSEQ
██████████	LBEO	STUDYID	DOMAIN	USUBJID	LB.LBSEQ
██████████	LBBR	STUDYID	DOMAIN	USUBJID	LB.LBSEQ

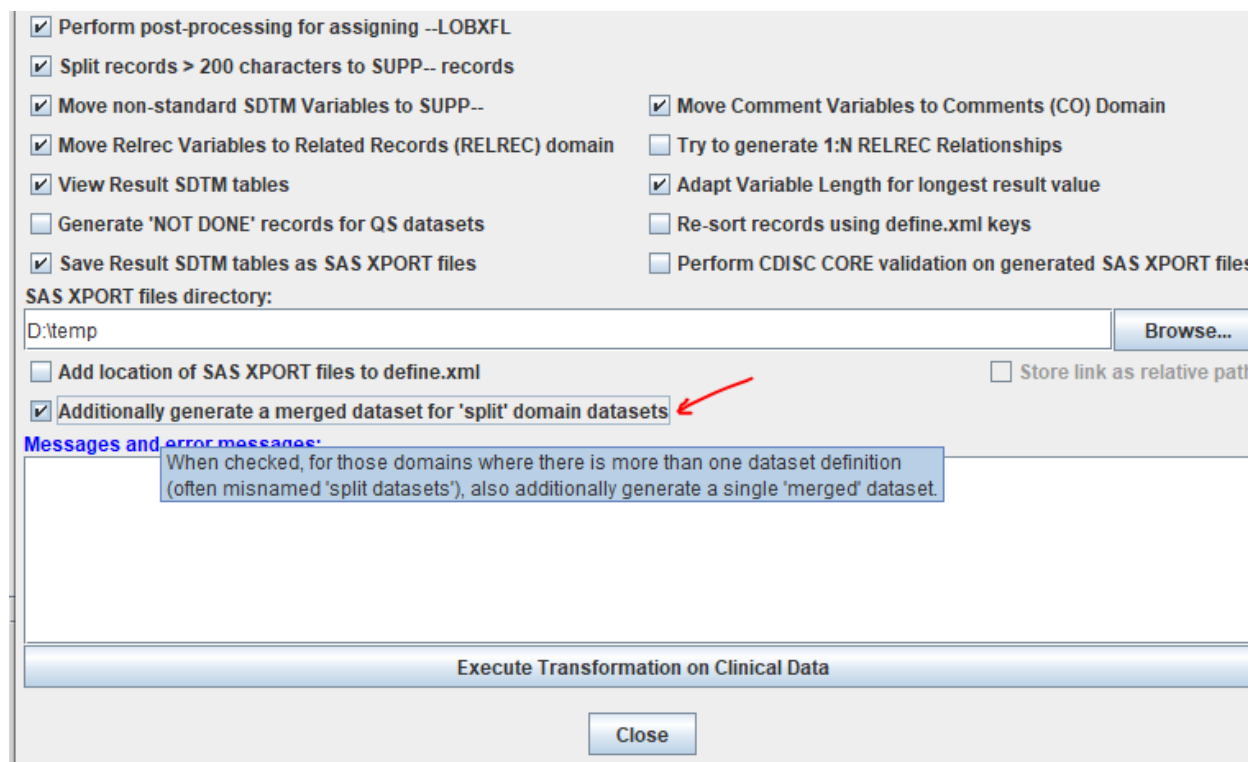
where it was decided to develop mappings and generate 4 instances of the LB (Laboratory) domain: one for urinalysis, one for all all blood tests (chemistry and hematology), one for coagulation tests, and one for breath alcohol tests.

Of course the choice of categories is arbitrary. For example, we could also have chosen to generate separate datasets for hematology and blood chemistry.

With the above setup, developing the mappings will be considerably easy, and lead to 4 submission datasets. In case SAS Transport is used: LBUR.xpt, LBBL.xpt, LBCO.xpt and LBBR.xpt. These will be considerably easier to analyze by the regulatory reviewers than would be the case with one huge LB dataset.

However, these authorities (or in case work is done as a service provider, the sponsor) will also want to see a single LB dataset.

In SDTM-ETL, this can be simply realized during execution of the mappings (when the XPT datasets are generated), by checking the checkbox "Additionally generate a merged dataset for 'split' domain datasets".



The screenshot shows the SAS SDTM-ETL configuration window. The 'SAS XPORT files directory' is set to 'D:\temp'. The checkbox 'Additionally generate a merged dataset for 'split' domain datasets' is checked and highlighted with a red arrow. A tooltip explains: 'When checked, for those domains where there is more than one dataset definition (often misnamed 'split datasets'), also additionally generate a single 'merged' dataset.' The window also contains a 'Close' button and a 'Execute Transformation on Clinical Data' button.

When it is checked, also a single dataset named "LB_merge.xpt" is generated in the output folder, containing the combined content of the separate LBUR, LBBL, etc. datasets. If wanted, this dataset can then of course be renamed, e.g. into LB.xpt.

Remark that for QS (Questionnaires), one will always generate distinct datasets, one for each questionnaire, and will never merge these into a single dataset.

Merging datasets that were generated separately

The above approach works very well when the different instances of the same domain have been developed together, or are loaded by merging (using the menu "File - Load Study define.xml" followed by "I want to merge with the existing define.xml").

There may be circumstances however, that also this approach does not work.

One such case is that we have two (or more) types of data in SUPPxx:

- SUPPxx dataset with "non-standard variables" automatic split off during mapping execution, and maybe containing values beyond the 200 character limitation.

- custom SUPPxx datasets (derived from the "SUPPQUAL" in the template)

Another example is the case that we have generated different CO (Comments) datasets, by automatic split of ("Comment variable in the dataset definition" generated by "Insert - New SDTM Variable for Comment" - see the tutorial "[Auto-generation of comments and putting them in the Comments \(CO\) domain](#)").

In this tutorial we will take the somewhat more complex example of a SUPPLB dataset that has the "Reason for Event" for each measurement (**one record per measurement** when the field on the CRF was filled), to be combined with a SUPPLB dataset that contains the "abnormality clinical significant" values, with **one record per visit** per group of records (i.e. CRF page). So, in the first SUPPLB dataset, the identifying variable (IDVAR) will be LBSEQ (Sequence Number), whereas in the second, it will be VISITNUM (Visit Number).

For this new feature, we made a separate application, named "XML2SASDatasetMerger". It does not merge SAS-XPT datasets directly (that should be done in SAS), but allows to use an XML representation of the output of SDTM-ETL (even when generating XPT files) and merge these. So no (expensive) SAS license is necessary.

Suppose we have already loaded our mappings for LB:

RELSUB	STUDYID	USUBJID	RELSUB.POOLID	RELSUB.RSUB...	RELSUB
OI	STUDYID	DOMAIN	OI.NHOID	OI.OISEQ	OI.OIPAR
██████████ DM	STUDYID	DOMAIN	USUBJID	SUBJID	DM.RFST
██████████ GLOBAL	STARTREF				
██████████ LBUR	STUDYID	DOMAIN	USUBJID	LB.LBSEQ	LB.LBGR
██████████ LBBL	STUDYID	DOMAIN	USUBJID	LB.LBSEQ	LB.LBGR
██████████ LBEO	STUDYID	DOMAIN	USUBJID	LB.LBSEQ	LB.LBGR
██████████ LBEBR	STUDYID	DOMAIN	USUBJID	LB.LBSEQ	LB.LBGR


We also loaded the DM dataset, in order to derive the "Last Observation Before First Exposure Flag" (LBLOBXFL) values, as that requires the value of DM.RFXSTDTC (datetime of first exposure).

We also see that an "non-standard variable" (NSV) REASEV" (Reason for Event) has been added to each of the instances:

			define.xml information: SDTM Name: REASEV OID: LB.REASEV Mandatory: No OrderNumber: 52 Role: SUPPQUAL Data type: text Length: 71 Description: Reason for Event Origin: Collected - Source: Investigator CRF
LB.LBRFTDTC	LB.LBEVINTX	LB.REASEV	
LB.LBRFTDTC	LB.LBEVINTX	LB.REASEV	
LB.LBRFTDTC	LB.LBEVINTX	LB.REASEV	
LB.LBRFTDTC	LB.LBEVINTX	LB.REASEV	

As it is an NSV, at least for regulatory submissions, it must be "banned" to SUPPLB.

When executing the mappings (using the menu "Transform - Generate Transformation (XSLT) Code for SAS-XPT", and following the wizard, this leads to the dialog:

Save output XML to file 

Perform post-processing for assigning --LOBXFL

Split records > 200 characters to SUPP-- records

Move non-standard SDTM Variables to SUPP--

Move Comment Variables to Comments (CO) Domain

Move Relrec Variables to Related Records (RELREC) domain

Try to generate 1:N RELREC Relationships

View Result SDTM tables

Adapt Variable Length for longest result value

Generate 'NOT DONE' records for QS datasets

Re-sort records using define.xml keys

Save Result SDTM tables as SAS XPORT files

Perform CDISC CORE validation on generated SAS XPORT files


SAS XPORT files directory:

Add location of SAS XPORT files to define.xml Store link as relative path

Additionally generate a merged dataset for 'split' domain datasets

Messages and error messages:

We see there is a checkbox and field "Save output XML to file:"
 Though we are generating SAS-XPT files, we can also save all the datasets into a single XML file, which essentially is a [CDISC Dataset-XML](#) file. For example:

Save output XML to file 

Perform post-processing for assigning --LOBXFL

Split records > 200 characters to SUPP-- records

Move non-standard SDTM Variables to SUPP--

Move Comment Variables to Comments (CO) Domain

Move Relrec Variables to Related Records (RELREC) domain

Try to generate 1:N RELREC Relationships

View Result SDTM tables

Adapt Variable Length for longest result value

Generate 'NOT DONE' records for QS datasets

Re-sort records using define.xml keys

Save Result SDTM tables as SAS XPORT files

Perform CDISC CORE validation on generated SAS XPORT files

SAS XPORT files directory:

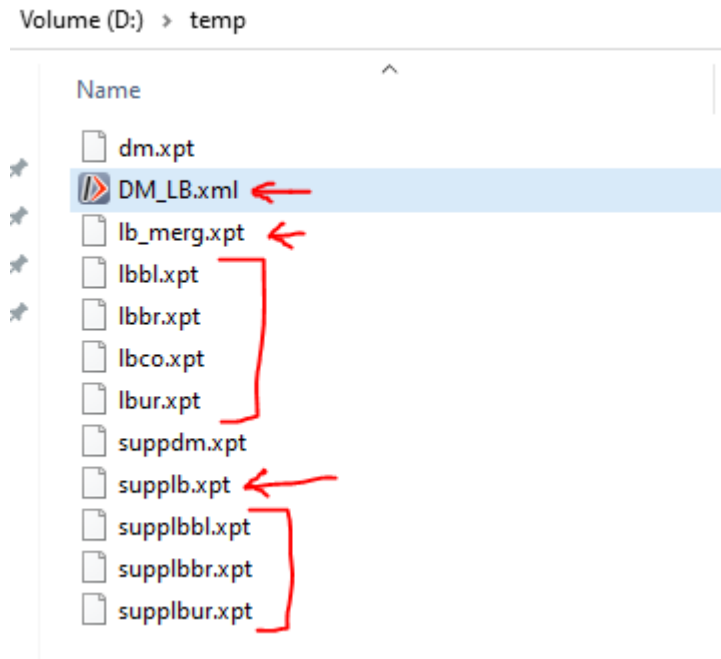
Add location of SAS XPORT files to define.xml Store link as relative path

Additionally generate a merged dataset for 'split' domain datasets

Messages and error messages:

We also check the checkbox "Additionally generate a merged dataset for 'split' domain datasets".

After execution, we find the following datasets:



We find the four (distinct) LB datasets: lbbl.xpt (blood tests), lbbr.xpt (breath alcohol tests), lbco.xpt (coagulation tests) and lbur.xpt (urinalysis), and the respective SUPPLB datasets supplbbl.xpt, supplbbr.xpt and supplbur.xpt. There is no supplbco.xpt dataset, as there were no "coagulation" records for which there was a value "Reason for Event".

For example, for supplbbl.xpt, we find:

	STUDYID	RDOMAIN	USUBJID	IDVAR	IDVARVAL	QNAM	QLABEL	QVAL	QORIG	QEVAL
1	[REDACTED]	LB	1006	LBSEQ	167	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
2	[REDACTED]	LB	1006	LBSEQ	168	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
3	[REDACTED]	LB	1006	LBSEQ	169	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
4	[REDACTED]	LB	1006	LBSEQ	170	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
5	[REDACTED]	LB	1006	LBSEQ	171	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
6	[REDACTED]	LB	1006	LBSEQ	172	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
7	[REDACTED]	LB	1006	LBSEQ	173	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
8	[REDACTED]	LB	1006	LBSEQ	174	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
9	[REDACTED]	LB	1006	LBSEQ	175	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
10	[REDACTED]	LB	1006	LBSEQ	176	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
11	[REDACTED]	LB	1006	LBSEQ	186	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
12	[REDACTED]	LB	1006	LBSEQ	187	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
13	[REDACTED]	LB	1006	LBSEQ	188	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
14	[REDACTED]	LB	1006	LBSEQ	189	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
15	[REDACTED]	LB	1006	LBSEQ	190	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	

also showing that each record is related to a single measurement, identified by LBSEQ, in the LB, and in the LBBL dataset.

We also see that a "merged" dataset named "lb_merg.xpt" has been created. This can be renamed to "lb.xpt" if wanted. And there is of course also a "merged" supplb.xpt dataset.

The additional XML dataset "DM_LB.xml" is a Dataset-XML file containing all the records of all the generated datasets, but then in XML format.

As once can see, it's file size is a bit higher than the sum of the file sizes of the XPT files, but not much more, which against contradicts the prejudice of the FDA that Dataset-XML files are much larger in size than XPT files².

The XML file "DM_LB.xml" contains all the records for DM, SUPPDM, LBBL, LBBR, LBCO, LBUR, SUPPLBBL, SUPPLBBR, SUPPLBUR, for example:

```
189859 <ItemGroupData ItemGroupOID="[REDACTED] SUPPDM" TransactionType="Insert">
189860   <ItemData ItemOID="STUDYID" Value="[REDACTED]" />
189861   <ItemData ItemOID="RDOMAIN" Value="DM" />
189862   <ItemData ItemOID="USUBJID" Value="4102" />
189863   <ItemData ItemOID="IDVAR" Value="" />
189864   <ItemData ItemOID="IDVARVAL" Value="" />
189865   <ItemData ItemOID="QNAM" Value="SCRNUM" />
189866   <ItemData ItemOID="QLABEL" Value="Screening Number" />
189867   <ItemData ItemOID="QVAL" Value="75" />
189868   <ItemData ItemOID="QORIG" Value="COLLECTED" />
189869   <ItemData ItemOID="QEVAL" Value="" />
189870 </ItemGroupData>
189871 <ItemGroupData ItemGroupOID="[REDACTED] SUPPLBUR" TransactionType="Insert">
189872   <ItemData ItemOID="STUDYID" Value="[REDACTED]" />
189873   <ItemData ItemOID="RDOMAIN" Value="LB" />
189874   <ItemData ItemOID="USUBJID" Value="4005" />
189875   <ItemData ItemOID="IDVAR" Value="LBSEQ" />
189876   <ItemData ItemOID="IDVARVAL" Value="53" />
189877   <ItemData ItemOID="QNAM" Value="REASEV" />
189878   <ItemData ItemOID="QLABEL" Value="Reason for Event" />
189879   <ItemData ItemOID="QVAL" Value="Former measurement: Reserve subject" />
189880   <ItemData ItemOID="QORIG" Value="COLLECTED" />
189881   <ItemData ItemOID="QEVAL" Value="" />
189882 </ItemGroupData>
```

each record having an identifier XXX:yyyy where XXX is the StudyID (hidden here) and yyy is the dataset identifier.

However, we also needed to generate a separate SUPPLBCS dataset ("CS" standing for "Clinically Significant"), as we have a field like "(Clinically significant) XXX lab value abnormalities?" which is only asked once per visit, and where XXX can e.g. be "glucose", "coagulation", "hematology".

As these are not directly related to a single measurement value, we choose to generate the SUPPLBCS dataset with the identifying variable IDVAR=VISITNUM³.

We generated this SUPPLBCS dataset separately in SDTM-ETL, using the menu "Insert - Domain specific SUPPQUAL", and then used QNAM as the variable to iterate over, i.e. a record is generated for each entry in the ODM dataset for which there is a mapping defined in QNAM:

² Also, XML files can be compressed to less than 5% of their normal size, and do not need to be de-compressed in order to be ready by modern software. This was however fully ignored by FDA during the Dataset-XML pilot. Essentially, the file sizes were just an excuse not having to change anything ...

³ Another solution would have been to add a variable LBGRPID in each of the LB instances, and reference that.

RELREC	STUDYID	RDOMAIN	USUBJID	IDVAR	IDVARVAL	RELTYPE	RELID	
SUPPQUAL	STUDYID	RDOMAIN	USUBJID	IDVAR	IDVARVAL	SUPPQUAL.QNAM	SUPPQUAL.QLABEL	ST
RELSUB	STUDYID	USUBJID	RELSUB.POOLID	RELSUB.RSUB...	RELSUB.SREL			
OI	STUDYID	DOMAIN	OI.NHOID	OI.OISEQ	OI.OIPARMCD	OI.OIPARM	OI.OIVAL	
■■■■■ SUPPLBCS	STUDYID	RDOMAIN	USUBJID	IDVAR	IDVARVAL	SUPPQUAL.QNAM	SUPPQUAL.QLABEL	ST

Also here, when we generate the XPT datasets, we also want to save a copy in XML format:

Save output XML to file

D:\temp\SUPPLBCS.xml Browse...

Perform post-processing for assigning --LOBXFL

Split records > 200 characters to SUPP-- records

Move non-standard SDTM Variables to SUPP--

Move Comment Variables to Comments (CO) Domain

Move Relrec Variables to Related Records (RELREC) domain

Try to generate 1:N RELREC Relationships

View Result SDTM tables

Adapt Variable Length for longest result value

Generate 'NOT DONE' records for QS datasets

Re-sort records using define.xml keys

Save Result SDTM tables as SAS XPORT files

Perform CDISC CORE validation on generated SAS XPORT files

SAS XPORT files directory:

D:\temp Browse...

Add location of SAS XPORT files to define.xml Store link as relative path

Additionally generate a merged dataset for 'split' domain datasets

[Messages and error messages:](#)

and the generated XPT (supplbcs.xpt) file looks like:

SAS Universal Viewer - [supplbcs.xpt]

File Tools Window Help

Address

Library Properties SUPPLBCS

Freeze Hide Show... Format Filter... Font... Find

Table View

	STUDYID	RDOMAIN	USUBJID	IDVAR	IDVARVAL	QNAM	QLABEL	QVAL	QORIG
1	■■■■■	LB	1001	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF
2	■■■■■	LB	1001	VISITNUM	0	GLUCLSIG	(Clin. sign.) glucose lab abnormalities	No	CRF
3	■■■■■	LB	1001	VISITNUM	0	COACLSIG	(Clin. sign.) coagul. lab abnormalities	No	CRF
4	■■■■■	LB	1001	VISITNUM	0	HEMCLSIG	(Clin. sign.) hematomol. lab abnormalities	No	CRF
5	■■■■■	LB	1002	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF
6	■■■■■	LB	1002	VISITNUM	0	GLUCLSIG	(Clin. sign.) glucose lab abnormalities	No	CRF
7	■■■■■	LB	1002	VISITNUM	0	COACLSIG	(Clin. sign.) coagul. lab abnormalities	No	CRF
8	■■■■■	LB	1002	VISITNUM	0	HEMCLSIG	(Clin. sign.) hematomol. lab abnormalities	No	CRF
9	■■■■■	LB	1003	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF
10	■■■■■	LB	1003	VISITNUM	0	GLUCLSIG	(Clin. sign.) glucose lab abnormalities	No	CRF
11	■■■■■	LB	1003	VISITNUM	0	COACLSIG	(Clin. sign.) coagul. lab abnormalities	No	CRF
12	■■■■■	LB	1003	VISITNUM	0	HEMCLSIG	(Clin. sign.) hematomol. lab abnormalities	No	CRF
13	■■■■■	LB	1004	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF
14	■■■■■	LB	1004	VISITNUM	0	GLUCLSIG	(Clin. sign.) glucose lab abnormalities	No	CRF

The corresponding XML file SUPPLBCS.xml has e.g.:


```

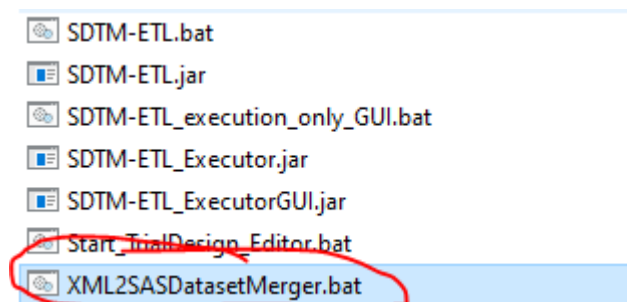
6 <ItemGroupData ItemGroupOID="████████ SUPPLBCS" TransactionType="Insert">
7   <ItemData ItemOID="STUDYID" Value="████████" />
8   <ItemData ItemOID="RDOMAIN" Value="LB"/>
9   <ItemData ItemOID="USUBJID" Value="1001"/>
10  <ItemData ItemOID="IDVAR" Value="VISITNUM"/>
11  <ItemData ItemOID="IDVARVAL" Value="0"/>
12  <ItemData ItemOID="SUPPQUAL.QNAM" Value="ABCLSIG"/>
13  <ItemData ItemOID="SUPPQUAL.QLABEL" Value="(Clin. sign.) urinal. dipstick abnormal."/>
14  <ItemData ItemOID="SUPPQUAL.QVAL" Value="No"/>
15  <ItemData ItemOID="SUPPQUAL.QORIG" Value="CRF"/>
16  <ItemData ItemOID="SUPPQUAL.QEVAL" IsNull="Yes"/>
17 </ItemGroupData>
18 <ItemGroupData ItemGroupOID="████████ SUPPLBCS" TransactionType="Insert">
19   <ItemData ItemOID="STUDYID" Value="████████" />
20   <ItemData ItemOID="RDOMAIN" Value="LB"/>
21   <ItemData ItemOID="USUBJID" Value="1001"/>
22   <ItemData ItemOID="IDVAR" Value="VISITNUM"/>
23   <ItemData ItemOID="IDVARVAL" Value="0"/>
24   <ItemData ItemOID="SUPPQUAL.QNAM" Value="GLUCLSIG"/>
25   <ItemData ItemOID="SUPPQUAL.QLABEL" Value="(Clin. sign.) glucose lab abnormalities"/>
26   <ItemData ItemOID="SUPPQUAL.QVAL" Value="No"/>
27   <ItemData ItemOID="SUPPQUAL.QORIG" Value="CRF"/>
28   <ItemData ItemOID="SUPPQUAL.QEVAL" IsNull="Yes"/>
29 </ItemGroupData>

```

Suppose now we want to generate a single SUPPLB-XPT file that contains as well the "REASEV" Reason for Event (REASEV) as the "clinical significance of xxx abnormalities". This means that we need to merge the SUPPLB records from both files.

We will not try to merge the XPT datasets directly (this would require an over-expensive SAS license), but we will use our new program "XML2SASDatasetMerger" which comes with the new SDTM-ETL v.4.3 distribution.

It can be started from the distribution folder by double-clicking (when using Windows) "XML2SASDatasetMerger.bat":



When executed, a new window is opened:

As soon as we add the location of the define.xml file, the combobox "Dataset OID / Name" presents a list with all the dataset definitions found in that define.xml:

The maximal length for each of the variables will be determined from the data in the files themselves, so that the generated XPT file will already have been optimized for file size. Please see the tutorial 'Merging similar datasets' for more details.

Dataset-XML File 1:	D:\temp\DM_LB.xml	Browse
Dataset-XML File 2:	D:\temp\SUPPLBCS.xml	Browse
Define-XML File:	D:\temp\SUPPLBCS_define.xml	Browse
Dataset-OID / Name:	CO	
Output SAS-XPT File:	CO	Browse
	DM	
	SE	
OID	SM	Data type
STUDYID	SV	Char
DOMAIN	AG	Char
RDOMAIN	CM	Char
USUBJID	EX	Char
CO.COSEQ		Num
CO.IDVAR		Char

We scroll down, and look for SUPPLBxx, and select it:

The maximal length for each of the variables will be determined from the data in the files themselves, so that the generated XPT file will already have been optimized for file size. Please see the tutorial 'Merging similar datasets' for more details.

Dataset-XML File 1:	D:\temp\DM_LB.xml	Browse
Dataset-XML File 2:	D:\temp\SUPPLBCS.xml	Browse
Define-XML File:	D:\temp\SUPPLBCS_define.xml	Browse
Dataset-OID / Name:	CO	
Output SAS-XPT File:	TM	Browse
	TI	
	TS	
OID	RELREC	Data type
STUDYID	SUPPQUAL	Char
DOMAIN	RELSUB	Char
RDOMAIN	COI	Char
USUBJID	SUPPLBCS	Char
CO.COSEQ		Num
CO.IDVAR		Char
CO.IDVARVAL	IDVARVAL	Identifying Variable Value
		Char

with "SUPPLBCS" prefixed by the StudyID (hidden here). Immediately, the table lower down updates to:

The maximal length for each of the variables will be determined from the data in the files themselves, so that the generated XPT file will already have been optimized for file size. Please see the tutorial 'Merging similar datasets' for more details.

Dataset-XML File 1:	D:\temp\DM_LB.xml	Browse
Dataset-XML File 2:	D:\temp\SUPPLBCS.xml	Browse
Define-XML File:	D:\temp\SUPPLBCS_define.xml	Browse
Dataset-OID / Name:	████████ SUPPLBCS	▼
Output SAS-XPT File:		Browse

OID	Name	Label	Data type
STUDYID	STUDYID	Study Identifier	Char
RDOMAIN	RDOMAIN	Related Domain Abbreviation	Char
USUBJID	USUBJID	Unique Subject Identifier	Char
IDVAR	IDVAR	Identifying Variable	Char
IDVARVAL	IDVARVAL	Identifying Variable Value	Char
SUPPQUAL.QNAM	QNAM	Qualifier Variable Name	Char
SUPPQUAL.QLABEL	QLABEL	Qualifier Variable Label	Char
SUPPQUAL.QVAL	QVAL	Data Value	Char
SUPPQUAL.QORIG	QORIG	Origin	Char
SUPPQUAL.QEVAL	QEVAL	Evaluator	Char

showing all the variables and their properties. This table will then be used to assign the "labels" in the XPT dataset.

We then also add the name and location of the XPT dataset we want to produce, e.g.:

Output SAS-XPT File:	D:\temp\supplb_overall.xpt	Browse
----------------------	----------------------------	--------

OID	Name	Label	Data type
STUDYID	STUDYID	Study Identifier	Char
RDOMAIN	RDOMAIN	Related Domain Abbreviation	Char
USUBJID	USUBJID	Unique Subject Identifier	Char
IDVAR	IDVAR	Identifying Variable	Char
IDVARVAL	IDVARVAL	Identifying Variable Value	Char
SUPPQUAL.QNAM	QNAM	Qualifier Variable Name	Char
SUPPQUAL.QLABEL	QLABEL	Qualifier Variable Label	Char
SUPPQUAL.QVAL	QVAL	Data Value	Char
SUPPQUAL.QORIG	QORIG	Origin	Char
SUPPQUAL.QEVAL	QEVAL	Evaluator	Char

The next step is a very important, as it allows to steer the selection of the records that will go into the output XPT dataset.

In the lower field we find:

SUPPQUAL.QORIG	QORIG	Origin	Char
SUPPQUAL.QEVAL	QEVAL	Evaluator	Char

'Starts With' OID of records to be retained and merged:

Start Merging

with XXX:SUPPLBCS just being a first proposal ...
and XXX being the StudyID (hidden here).

As we have seen before, all our "REASEV" records in the Dataset-XML file "DM_LB.xml" have the identifier (OID) either being "XXX:SUPPLBBL", "XXX:SUPPLBBR", or "XXX:SUPPLBUR",
and the one for the "Significance" records in dataset "SUPPLBCS.xml" have the identifier "XXX:SUPPLBCS".

The "Starts With" field now allows us to select which records go into the XPT dataset that we want to generate. If we leave it being "XXX:SUPPLBCS", we will get a SUPPLB dataset that only has "Clinical Significant" records. So, we change it into:

'Starts With' OID of records to be retained and merged: <input type="text" value="XXXXXXXXX SUPPLB"/>	All records (ItemGroupData) in both files for which the OID <u>starts with</u> the given string (case sensitive), will be retained and copied into the output SAS-XPT file.
Start Merging	You may edit the value in order to be less / more selective in which records are retained

stating that all records for which the identifier "starts with" "XXX:SUPPLB" will be selected, which means all "REASEV" records in file "DM_LB.xml" as well as all "Significance" records in the file "SUPPLBCS.xml".

When then clicking "Start Merging", the software starts the following process:

- it first analyzes the length for each value of each variable from the XML files, in order to determine the maximal value length, that SAS-XPT requires⁴.
- it then selects the records from each dataset and combines them into a single XML file (if you want to see it, it is named "temp_merge.xml" and is in the "temp" folder of your SDTM-ETL distribution folder).
- it then takes this combined XML dataset, and with the given labels and determined maximal lengths for the variables, and generates the SAS-XPT dataset and writes it in the output file.

The result in this case e.g. is:

⁴ "Maximal length" of the variable value would be completely irrelevant when working with modern data formats like XML, JSON, YAML, ...

SAS Universal Viewer - [supplb_overall.xpt]

File Tools Window Help

Address

Library Properties SUPPLB

Freeze Hide Show... Format Filter... Font... Find

Table View

	STUDYID	RDOMAIN	USUBJID	IDVAR	IDVARVAL	QNAM	QLABEL	QVAL	QORIG	QEV
28	██████████	LB	1006	LBSEQ	193	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
29	██████████	LB	1006	LBSEQ	194	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
30	██████████	LB	1006	LBSEQ	195	REASEV	Reason for Event	Extra measurement: Safety	COLLECTED	
31	██████████	LB	2003	LBSEQ	167	REASEV	Reason for Event	██████████redose.	COLLECTED	
32	██████████	LB	2003	LBSEQ	168	REASEV	Reason for Event	██████████redose.	COLLECTED	
33	██████████	LB	2003	LBSEQ	169	REASEV	Reason for Event	██████████redose.	COLLECTED	
34	██████████	LB	2003	LBSEQ	170	REASEV	Reason for Event	██████████redose.	COLLECTED	
35	██████████	LB	2003	LBSEQ	171	REASEV	Reason for Event	██████████redose.	COLLECTED	
36	██████████	LB	2003	LBSEQ	172	REASEV	Reason for Event	██████████redose.	COLLECTED	
37	██████████	LB	2003	LBSEQ	173	REASEV	Reason for Event	██████████redose.	COLLECTED	
38	██████████	LB	2003	LBSEQ	174	REASEV	Reason for Event	██████████redose.	COLLECTED	
39	██████████	LB	2003	LBSEQ	175	REASEV	Reason for Event	██████████redose.	COLLECTED	
40	██████████	LB	2003	LBSEQ	176	REASEV	Reason for Event	██████████redose.	COLLECTED	
41	██████████	LB	4003	LBSEQ	202	REASEV	Reason for Event	Fomer measurement: Reserve subject Fomer measurement: Reserve subject	COLLECTED	
42	██████████	LB	4005	LBSEQ	218	REASEV	Reason for Event	Fomer measurement: Reserve subject Fomer measurement: Reserve subject	COLLECTED	
43	██████████	LB	1001	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF	
44	██████████	LB	1001	VISITNUM	0	GLUCLSIG	(Clin. sign.) glucose lab abnormalities	No	CRF	
45	██████████	LB	1001	VISITNUM	0	COACLSIG	(Clin. sign.) coagul. lab abnormalities	No	CRF	
46	██████████	LB	1001	VISITNUM	0	HEMCLSIG	(Clin. sign.) hematom. lab abnormalities	No	CRF	
47	██████████	LB	1002	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF	
48	██████████	LB	1002	VISITNUM	0	GLUCLSIG	(Clin. sign.) glucose lab abnormalities	No	CRF	
49	██████████	LB	1002	VISITNUM	0	COACLSIG	(Clin. sign.) coagul. lab abnormalities	No	CRF	
50	██████████	LB	1002	VISITNUM	0	HEMCLSIG	(Clin. sign.) hematom. lab abnormalities	No	CRF	
51	██████████	LB	1003	VISITNUM	0	ABCLSIG	(Clin. sign.) urinal. dipstick abnormal.	No	CRF	

And as one can see from the properties, the combined SUPPLB file (supplb_overall.xpt) has been optimized for "minimizing" the file size:

SAS Universal Viewer - [supplb_overall.xpt]

File Tools Window Help

Address

Library Properties SUPPLB

Name	Obs	Vars	#	Variable	Type	Length	Format	Infomat	Label
SUPPLB	143	10	1	STUDYID	Character	8			Study Identifier
			2	RDOMAIN	Character	2			Related Domain Abbreviation
			3	USUBJID	Character	4			Unique Subject Identifier
			4	IDVAR	Character	8			Identifying Variable
			5	IDVARVAL	Character	3			Identifying Variable Value
			6	QNAM	Character	8			Qualifier Variable Name
			7	QLABEL	Character	40			Qualifier Variable Label
			8	QVAL	Character	71			Data Value
			9	QORIG	Character	9			Origin
			10	QEV	Character	1			Evaluator

We have demonstrated the XML2SASDatasetMerger software for the more complicated case of merging two different types of SUPPLB datasets. The same approach can however be used for merging e.g. different instances of CO (Comments) that have been generated separately for different domains in different datasets.

The only current limitation is that the datasets to be merged have the same (number of) variables. For example, it will usually not allow to e.g. merge an LB dataset with an MB dataset.

Conclusions

The by FDA and other regulatory authorities mandated (but outdated) SAS-XPT format makes life unnecessary complicated when generating CDISC SDTM and SEND datasets.

Essentially, "splitting" and "merging" datasets should never have to be done, and if so anyway, it should be easy to accomplish. It is the SAS-XPT format that makes it complicated: when using XML (e.g. Dataset-XML) or JSON (e.g. Dataset-JSON), this is much easier to accomplish.

This tutorial demonstrated two ways of merging datasets for the case of XPT as the output format. For the simple case that several instances of the same domain are present in the same "working" define.xml, one can set the *"Additionally generate a merged dataset for 'split' domain datasets"* checkbox.

In case datasets to be merged have been created separately, one can use the new utility tool (as of SDTM-ETL v.4.3) named "XML2SASDatasetMerger" that comes as a separate program.

We expect that FDA will give green light for Dataset-JSON submissions by the end of this year. As soon as that is clear, we will also provide functionalities for merging Dataset-JSON files. Developing these will however be much more straightforward than for SAS-XPT.